

11-15-00

Attorney Docket No. YOR920000401US1
**IN THE UNITED STATES
PATENT AND TRADEMARK OFFICE**

11/14/00
35c93 U.S. PTO

PTO
U.S. PTO
09/13/00
11/14/00

PATENT APPLICATION

"Express Mail" Label No.: EL659921972US
I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 C.F.R. §1.10 on the date indicated below and is addressed to: Assistant Commissioner for Patents, Washington, D.C. 20231.

APPLICANT(S): Youssef Drissi and Ricardo Vilalta

TITLE: **METHODS AND
APPARATUS FOR
GENERATING A DATA
CLASSIFICATION MODEL
USING AN ADAPTIVE
LEARNING ALGORITHM**

Date of Deposit: November 14, 2000

Signature: Linda M. Schubert

**ASSISTANT COMMISSIONER FOR PATENTS
WASHINGTON, D.C. 20231**

SIR:

Enclosed are the following papers relating to the above-named application for patent:

- Application - 11 pages of Specification, 6 pages of Claims and 1 page of Abstract
- 9 Sheets of informal drawing(s)
- Declaration
- 1 Assignment with Cover Sheet
- Associate Power of Attorney

CLAIMS AS FILED				
	NO. FILED	NO. EXTRA	RATE	CALCULATIONS
Total Claims	23 -20 =	3	x \$18 =	\$ 54
Independent Claims	7 - 3 =	4	x \$80 =	\$ 320
Multiple Dependent Claim(s), if applicable			\$260 =	\$ 0
Basic Fee				\$ 710
TOTAL FEE:				\$1,084

Please file the application and charge **IBM Corporation's Deposit Account No. 50-0510** the amount of \$1,084, to cover the filing fee. In the event of non-payment or improper payment of a required fee, the Commissioner is authorized to charge or to credit **Deposit Account No. 50-0510** as required to correct the error. *Duplicate copies of this letter are enclosed.*

Please address all correspondence to: **Kevin M. Mason, Ryan, Mason & Lewis, LLP, 1300 Post Road, Suite 205, Fairfield, CT 06430.** Telephone calls should be made to the under-signed attorney at (203) 255-6560.

Respectfully,

Kevin M. Mason
Kevin M. Mason
Reg. No. 36,597
Attorney for Applicant(s)

Date: November 14, 2000
Ryan, Mason & Lewis, LLP
1300 Post Road, Suite 205
Fairfield, CT 06430

09713342 11/14/00

**METHODS AND APPARATUS FOR GENERATING A DATA
CLASSIFICATION MODEL USING AN ADAPTIVE LEARNING ALGORITHM**

Cross Reference to Related Application

5 The present invention is related to United States Patent Application
entitled "Method and Apparatus for Generating a Data Classification Model Using
Interactive Adaptive Learning Algorithms," (Attorney Docket Number
YOR920000507US1), filed contemporaneously herewith, assigned to the assignee of the
present invention and incorporated by reference herein.

10 **Field of the Invention**

 The present invention relates generally to the fields of data mining or
machine learning and, more particularly, to methods and apparatus for generating data
classification models.

15 **Background of the Invention**

 Data classification techniques, often referred to as supervised learning,
attempt to find an approximation or hypothesis to a target concept that assigns objects
(such as processes or events) into different categories or classes. Data classification can
20 normally be divided into two phases, namely, a learning phase and a testing phase. The
learning phase applies a learning algorithm to training data. The training data is typically
comprised of descriptions of objects (a set of feature variables) together with the correct
classification for each object (the class variable).

 The goal of the learning phase is to find correlations between object
25 descriptions to learn how to classify the objects. The training data is used to construct
models in which the class variable may be predicted in a record in which the feature
variables are known but the class variable is unknown. Thus, the end result of the
learning phase is a model or hypothesis (e.g., a set of rules) that can be used to predict the

class of new objects. The testing phase uses the model derived in the training phase to predict the class of testing objects. The classifications made by the model is compared to the true object classes to estimate the accuracy of the model.

Numerous techniques are known for deriving the relationship between the
5 feature variables and the class variables, including, for example, Disjunctive Normal Form (DNF) Rules, decision trees, nearest neighbor, support vector machines (SVMs) and Bayesian classifiers, as described, for example, in R. Agrawal et al., "An Interval Classifier for Database Mining Applications," Proc. of the 18th VLDB Conference, Vancouver, British Columbia, Canada 1992; C. Apte et al., "RAMP: Rules Abstraction
10 for Modeling and Prediction," IBM Research Report RC 20271, June 1995; J.R. Quinlan, "Induction of Decision Trees," Machine Learning, Volume 1, Number 1, 1986; J. Shafer et al., "SPRINT: A Scaleable Parallel Classifier for Data Mining," Proc. of the 22d VLDB Conference, Bombay, India, 1996; M. Mehta et al., "SLIQ: A Fast Scaleable Classifier for Data Mining," Proceedings of the Fifth International Conference on
15 Extending Database Technology, Avignon, France, March, 1996, each incorporated by reference herein.

Data classifiers have a number of applications that automate the labeling of unknown objects. For example, astronomers are interested in automated ways to classify objects within the millions of existing images mapping the universe (e.g.,
20 differentiate stars from galaxies). Learning algorithms have been trained to recognize these objects in the training phase, and used to predict new objects in astronomical images. This automated classification process obviates manual labeling of thousands of currently available astronomical images.

While such learning algorithms derive the relationship between the feature
25 variables and the class variables, they generally produce the same output model given the same domain dataset. Generally, a learning algorithm encodes certain assumptions about the nature of the concept to learn, referred to as the bias of the learning algorithm. If the

assumptions are wrong, however, then the learning algorithm will not provide a good approximation of the target concept and the output model will exhibit low accuracy. Most research in the area of data classification has focused on producing increasingly more accurate models, which is impossible to attain on a universal basis over all possible domains. It is now well understood that increasing the quality of the output model on a certain group of domains will cause a decrease of quality on other groups of domains. See, for example, C. Schaffer, "A Conservation Law for Generalization Performance," Proc. of the Eleventh Int'l Conference on Machine Learning, 259-65, San Francisco, Morgan Kaufman (1994); and D. Wolpert, "The Lack of a Priori Distinctions Between Learning Algorithms and the Existence of a Priori Distinctions Between Learning Algorithms," Neural Computation, 8 (1996), each incorporated by reference herein.

While conventional learning algorithms produce sufficiently accurate models for many applications, they suffer from a number of limitations, which, if overcome, could greatly improve the performance of the data classification and regression systems that employ such models. Specifically, the learning algorithms of conventional data classification and regression systems are unable to adapt over time. In other words, once a model is generated by a learning algorithm, the model cannot be reconfigured based on experience. Thus, the conventional data classification and regression systems that employ such models are prone to repeating the same errors.

A need therefore exists for data classification and regression methods and apparatus that adapt a learning algorithm through experience. Another need exists for data classification and regression methods and apparatus that dynamically modify the assumptions of the learning algorithm to improve the assumptions embodied in the generated models and thereby improve the quality of the data classification and regression systems that employ such models. Yet another need exists for a learning method and apparatus that performs meta-learning to improve the assumptions or inductive bias in a model.

Summary of the Invention

Generally, a data classification method and apparatus are disclosed for labeling unknown objects. The disclosed data classification system employs a learning
5 algorithm that adapts through experience. The present invention classifies objects in domain datasets using data classification models having a corresponding bias and evaluates the performance of the data classification. The performance values for each domain dataset and corresponding model bias are processed to initially identify (and over
10 time modify) one or more rules of experience. The rules of experience are then subsequently used to generate a model for data classification. Each rule of experience specifies one or more characteristics for a domain dataset and a corresponding bias that should be utilized for a data classification model if the rule is satisfied.

Thus, the present invention dynamically modifies the assumptions (bias) of the learning algorithm to improve the assumptions embodied in the generated models
15 and thereby improve the quality of the data classification and regression systems that employ such models. Furthermore, since the rules of experience change dynamically, the learning process of the present invention will not necessarily output the same model when the same domain dataset is presented again. Furthermore, the disclosed self-adaptive learning process will become increasingly more accurate as the rules of experience are
20 accumulated over time.

A more complete understanding of the present invention, as well as further features and advantages of the present invention, will be obtained by reference to the following detailed description and drawings.

Brief Description of the Drawings

FIG. 1 is a schematic block diagram showing the architecture of an illustrative data classification system in accordance with the present invention;

FIG. 2 illustrates the operation of the data classification system;
FIG. 3 illustrates an exemplary table from the domain dataset of FIG. 1;
FIG. 4 illustrates an exemplary table from the performance dataset of FIG.
1;
5 FIG. 5 illustrates an exemplary table from the rules of experience table of
FIG. 1;
FIG. 6 is a flow chart describing the meta-feature generation process of
FIG. 1;
FIG. 7 is a flow chart describing the performance assessment process of
10 FIG. 1;
FIG. 8 is a flow chart describing the rules of experience generation process
of FIG. 1; and
FIG. 9 is a flow chart describing the self-adaptive learning process of FIG.
1 incorporating features of the present invention.

15

Detailed Description of Preferred Embodiments

FIG. 1 illustrates a data classification system 100 in accordance with the
present invention. The data classification system 100 may be embodied as a conventional
data classification system, such as the learning program described in J. R. Quinlan, C4.5:
20 Programs for Machine Learning. Morgan Kaufmann Publishers, Inc. Palo Alto, CA,
incorporated by reference herein, as modified in accordance with the features and
functions of the present invention to provide an adaptive learning algorithm.

FIG. 1 is a schematic block diagram showing the architecture of an
illustrative data classification system 100 in accordance with the present invention. The
25 data classification system 100 may be embodied as a general purpose computing system,
such as the general purpose computing system shown in FIG. 1. The data classification
system 100 includes a processor 110 and related memory, such as a data storage device

120, which may be distributed or local. The processor 110 may be embodied as a single processor, or a number of local or distributed processors operating in parallel. The data storage device 120 and/or a read only memory (ROM) are operable to store one or more instructions, which the processor 110 is operable to retrieve, interpret and execute. As shown in FIG. 1, the data classification system 100 optionally includes a connection to a computer network (not shown).

As shown in FIG. 1 and discussed further below in conjunction with FIGS. 3 through 5, the data storage device 120 preferably includes a domain dataset 300, a performance dataset 400 and a rules of experience table 500. Generally, the domain dataset 300 contains a record for each object and indicates the class associated with each object. The performance dataset 400 indicates the learning algorithm that produced the best model for each domain. The rules of experience table 500 identify a number of prioritized rules and their corresponding conditions, which if satisfied, provide a bias or assumption that should be employed when generating a model.

In addition, as discussed further below in conjunction with FIGS. 6 through 11, the data storage device 120 includes a meta-feature generation process 600, a performance assessment process 700, a rules of experience generation process 800 and a self-adaptive learning process 900. Generally, the meta-feature generation process 600 processes each domain dataset to represent the domain as a set of meta-features. The performance assessment process 700 evaluates the performance of a given model for a given domain dataset described by a set of meta-features and stores the results in the performance dataset 400. The rules of experience generation process 800 evaluates the performance dataset 400 in order to modify or extend the current rules in the rules of experience table 500. The self-adaptive learning process 900 identifies the best model for a given domain dataset 300, based on the current rules of experience table 500.

FIG. 2 provides a global view of the data classification system 100. As shown in FIG. 2, a domain dataset 300, discussed below in conjunction with FIG. 3,

serves as input to the system 100. The domain dataset 300 is applied to a self-adaptive learning process 900, discussed below in conjunction with FIG. 9, during step 220 and a meta-feature generation process 600, discussed below in conjunction with FIG. 6, during step 240. Generally, the self-adaptive learning process 900 produces an output model 250 that can be used to predict the class labels of future examples. For a detailed discussion of suitable models 250, see, for example, J.R. Quinlan, C4.5: Programs for Machine Learning. Morgan Kaufmann Publishers, Inc. Palo Alto, CA. (1994) (decision trees); Weiss, Sholom and Indurkha, Nitin, "Optimized Rule Induction", Intelligent Expert, Volume 8, Number 6, pp. 61-69, 1993 (rules); and L.R. Rivest, "Learning Decision Lists", Machine Learning, 2, 3, 229-246, (1987) (decision lists), each incorporated by reference herein.

The meta-feature generation process 600 executed during step 240 represents the domain dataset 300 as a set of meta-features. The performance of the output model 250 is assessed during step 260 by the performance assessment process 700, discussed below in conjunction with FIG. 7, and the performance assessment is recorded in the performance dataset 400. The performance assessment process 700 executed during step 260 evaluates how much the output model 250 can be improved.

As shown in FIG. 2, the self-adaptive learning process 900 receives the following information as inputs: (i) the domain dataset 300; (ii) the meta-feature description of the domain dataset 300; and (iii) the performance dataset 400. As discussed further below in conjunction with FIG. 9, the self-adaptive learning process 900 can use these inputs to modify the underlying assumptions embodied in a given model, such that, if the same dataset 300 were to be presented again to the self-adaptive learning process 900 a more accurate model would be produced.

DATABASES

FIG. 3 illustrates an exemplary table from the domain dataset 300 that includes training examples, each labeled with a specific class. As previously indicated,

the domain dataset 300 contains a record for each object and indicates the class associated with each object. The domain dataset 300 maintains a plurality of records, such as records 305 through 320, each associated with a different object. For each object, the domain dataset 300 indicates a number of features in fields 350 through 365, describing each object in the dataset. The last field 370 corresponds to the class assigned to each object. For example, if the domain dataset 300 were to correspond to astronomical images to be classified as either stars or galaxies, then each record 305-320 would correspond to a different object in the image, and each field 350-365 would correspond to a different feature such as the amount of luminosity, shape or size. The class field 370 would be populated with the label of "star" or "galaxy."

FIG. 4 illustrates an exemplary table from the performance dataset 400. As previously indicated, the performance dataset 400 indicates the performance for each model on a domain. The performance dataset 400 maintains a plurality of records, such as records 405 through 415, each associated with a different model. For each model, the performance dataset 400 identifies the domain on which the model was utilized in field 450, as well as the underlying bias embodied in the model in field 455 and the performance assessment in field 460. Each domain can be identified in field 450, for example, using a vector of meta-features characterizing each domain (as produced by the meta-feature generation process 600).

FIG. 5 illustrates an exemplary table from the rules of experience table 500. The rules of experience table 500 identifies a number of prioritized rules and their corresponding conditions, which if satisfied, provide a bias or assumption that should be employed when generating a model. As shown in FIG. 5, the rules of experience table 500 includes a plurality of records, such as records 505 through 515, each associated with a different experience rule. For each rule identified in field 550, the rules of experience table 500 identifies the corresponding conditions associated with the rule in field 560 and

the bias or assumption that should be employed in a model when the rule is satisfied in field 570.

PROCESSES

FIG. 6 is a flow chart describing the meta-feature generation process 600.

- 5 As previously indicated, the meta-feature generation process 600 processes each set of domain data to represent the domain as a set of meta-features. As shown in FIG. 6, the meta-feature generation process 600 initially processes the domain dataset 300 during step 610 to store the information in a table. Thereafter, the meta-feature generation process 600 extracts statistics from the dataset 300 during step 620 that are then used to
- 10 generate meta-features during step 630. For a discussion of the generation of meta-features that are particularly relevant to the meta-learning phase, including concept variation or average weighted distance meta-features, as well as additional well-known meta-features, see, for example, United States Patent Application Serial Number 09/629,086, filed July 31, 2000, entitled "Methods and Apparatus for Selecting a Data
- 15 Classification Model Using Meta-Learning," assigned to the assignee of the present invention and incorporated by reference herein.

- FIG. 7 is a flow chart describing the performance assessment process 700. The performance assessment process 700 evaluates the performance of a given model for a given domain dataset and stores the results in the performance dataset 400. The process
- 20 700 initially receives a model 250 during step 710 and assesses empirically the performance of the model 250. In other words, the model 250 is used to classify objects during step 710, for which the classification is already known, so that an objective measure of the model performance may be obtained. Typically, the performance assessment corresponds to the estimated accuracy of the model 250.

- 25 As shown in FIG. 7, the domain is then processed during step 715 by the meta-feature generation process 600, discussed above in conjunction with FIG. 6, to obtain a vector of meta-features characterizing the domain. Thereafter, a new entry is

created in the performance dataset 400 during step 720 using (i) the meta-feature description of the domain on which the model 250 was utilized, (ii) the underlying bias embodied in the model and (iii) the performance assessment determined during step 710.

FIG. 8 is a flow chart describing an exemplary rules of experience generation process 800 that evaluates the performance dataset 400 in order to modify or extend the current rules in the rules of experience table 500. As shown in FIG. 8, the rules of experience generation process 800 initially evaluates the performance dataset 400 during step 810 to identify correlations between various domains (described by a set of meta-features) and their corresponding best inductive bias (model).

Generally, the rules of experience generation process 800 employs a simple learning algorithm that receives a domain as input (in this case, the performance dataset 400) and produces as a result a model (in this case, the rule of experience 500). The difference lies in the nature of the domain. For a simple learning algorithm, the domain is a set of objects that belong to a real-world application, and where we wish to be able to predict the class of new objects. In the rules of experience generation process 800, each object contains the meta-features of a domain and the class of each object indicates the bias used to learn that domain. The rules of experience generation process 800 is thus a meta-learner that learns about the learning process itself. The mechanism behind it, however, is no different from a simple learning algorithm.

Based on the correlations identified during step 810, the current rules of experience are modified or extended during step 820 and recorded in the rules of experience table 500. For example, as shown in the exemplary rules of experience table 500 of FIG. 5, when models used a particular bias that partitioned the data in a specified manner, certain correlations were identified in various meta-features.

The modification or extension of the rules in the rules of experience table 500 will influence the future selection of models by the self-adaptive learning process 900, discussed below in conjunction with FIG. 9. Since the rules of experience change

dynamically, the learning process 900 of the present invention will not necessarily output the same model when the same domain dataset is presented again. Furthermore, the self-adaptive learning process 900 will become increasingly more accurate as the rules of experience table 500 grows larger.

5 FIG. 9 is a flow chart describing an exemplary self-adaptive learning process 900 that identifies the best model for a given domain dataset 300, based on the current rules of experience table 500. As shown in FIG. 9, the self-adaptive learning process 900 initially executes the meta-feature generation process 600, discussed above in conjunction with FIG. 6, during step 910 to provide a meta-feature description of the
10 current domain. During step 920, the self-adaptive learning process 900 sequentially compares the meta-feature description of the current domain to each of the rules in the rules of experience table 500 until a rule is satisfied. In this manner, the first satisfied rule provides the best bias to utilize for the current domain.

 If a rule is satisfied, then the corresponding bias is applied to generate the
15 model 250 during step 930. If, however, no rule in the rules of experience table 500 is satisfied for the current domain, then a default bias is retrieved during step 940 and the default bias is applied to generate the model 250 during step 930. Thereafter, program control terminates.

 It is to be understood that the embodiments and variations shown and
20 described herein are merely illustrative of the principles of this invention and that various modifications may be implemented by those skilled in the art without departing from the scope and spirit of the invention.

Claims

What is claimed is:

1. A method for classifying data, comprising the steps of:
5 classifying objects in a domain dataset using a data classification model,
said data classification model having a bias;
 evaluating the performance of said classifying step; and
 modifying said bias based on said performance evaluation.
- 10 2. The method of claim 1, wherein said steps of classifying and evaluating
are performed for a plurality of said domain datasets and wherein said method further
comprising the steps of recording a performance value for each combination of said
domain datasets and said bias.
- 15 3. The method of claim 2, further comprising the step of processing said
recorded performance values for each combination of said domain datasets and said bias
to generate one or more rules, each of said rules specifying one or more characteristics of
said domain datasets and a corresponding bias that should be utilized in one of said data
classification models.
- 20 4. The method of claim 3, further comprising the step of selecting a data
classification model for classifying a domain dataset by comparing characteristics of said
domain dataset to said rules.
- 25 5. The method of claim 1, wherein said domain dataset is represented using a
set of meta-features.

6. The method of claim 5, wherein said meta-features includes a concept variation meta-feature.

7. The method of claim 5, wherein said meta-features includes an average weighted distance meta-feature that measures the density of the distribution of said at least one domain dataset.

8. A method for classifying data, comprising the steps of:
classifying objects in a plurality of domain datasets using one of a number of data classification models, each of said data classification models having a corresponding bias;
evaluating the performance of each of said domain dataset classifications;
maintaining a performance value for each combination of said domain datasets and said bias;
processing said performance values for each combination of said domain datasets and said bias to generate one or more rules, each of said rules specifying one or more characteristics of said domain datasets and a corresponding bias that should be utilized in one of said data classification models; and
selecting a data classification model for classifying a domain dataset by comparing characteristics of said domain dataset to said rules.

9. The method of claim 8, further comprising the step of modifying at least one of said biases based on said performance evaluation.

10. The method of claim 8, wherein said domain dataset is represented using a set of meta-features.

11. The method of claim 10, wherein said meta-features includes a concept variation meta-feature.

12. The method of claim 10, wherein said meta-features includes an average weighted distance meta-feature that measures the density of the distribution of said at least one domain dataset.

13. A method for classifying data in a domain dataset, comprising:
applying an adaptive learning algorithm to said domain dataset to select a data classification model, said data classification model having a bias;
classifying objects in said domain dataset using said selected data classification model;
evaluating the performance of said classifying step;
maintaining an indication of said performance of said model for said domain dataset;
repeating said applying, classifying and evaluating steps for a plurality of said domain datasets; and
processing said performance values for each combination of said domain datasets and said bias to adjust one or more rules for subsequent data classification, each of said rules specifying one or more characteristics of said domain datasets and a corresponding bias that should be utilized in one of said data classification models.

14. The method of claim 13, further comprising the step of selecting a data classification model for classifying a domain dataset by comparing characteristics of said domain dataset to said rules.

15. The method of claim 13, further comprising the step of modifying at least one of said biases based on said performance evaluation.
16. A system for classifying data, comprising:
 - 5 a memory that stores computer-readable code; and
 - a processor operatively coupled to said memory, said processor configured to implement said computer-readable code, said computer-readable code configured to:
 - classify objects in a domain dataset using a data classification model, said data classification model having a bias;
 - 10 evaluate the performance of said classifying step; and
 - modify said bias based on said performance evaluation.
17. The system of claim 16, wherein said processor is further configured to classify said objects and evaluate said performance for a plurality of said domain datasets and wherein said processor records a performance value for each combination of said domain datasets and said bias.
18. The system of claim 17, wherein said processor is further configured to process said recorded performance values for each combination of said domain datasets and said bias to generate one or more rules, each of said rules specifying one or more characteristics of said domain datasets and a corresponding bias that should be utilized in one of said data classification models.
19. The system of claim 18, wherein said processor is further configured to select a data classification model for classifying a domain dataset by comparing characteristics of said domain dataset to said rules.

20. The system of claim 16, wherein said domain dataset is represented using a set of meta-features.

21. A system for classifying data, comprising:

5 a memory that stores computer-readable code; and
a processor operatively coupled to said memory, said processor configured to implement said computer-readable code, said computer-readable code configured to:

10 classify objects in a plurality of domain datasets using one of a number of data classification models, each of said data classification models having a corresponding bias;

evaluate the performance of each of said domain dataset classifications;

maintaining a performance value for each combination of said domain datasets and said bias;

15 process said performance values for each combination of said domain datasets and said bias to generate one or more rules, each of said rules specifying one or more characteristics of said domain datasets and a corresponding bias that should be utilized in one of said data classification models; and

select a data classification model for classifying a domain dataset by comparing characteristics of said domain dataset to said rules.

20

22. An article of manufacture for classifying data, comprising:

a computer readable medium having computer readable code means embodied thereon, said computer readable program code means comprising:

25 a step to classify objects in a domain dataset using a data classification model, said data classification model having a bias;

a step to evaluate the performance of said classifying step; and

a step to modify said bias based on said performance evaluation.

23. An article of manufacture for classifying data, comprising:

a computer readable medium having computer readable code means embodied thereon, said computer readable program code means comprising:

5 a step to classify objects in a plurality of domain datasets using one of a number of data classification models, each of said data classification models having a corresponding bias;

a step to evaluate the performance of each of said domain dataset classifications;

10 a step to maintaining a performance value for each combination of said domain datasets and said bias;

a step to process said performance values for each combination of said domain datasets and said bias to generate one or more rules, each of said rules specifying one or more characteristics of said domain datasets and a corresponding bias that should
15 be utilized in one of said data classification models; and

a step to select a data classification model for classifying a domain dataset by comparing characteristics of said domain dataset to said rules.

**METHODS AND APPARATUS FOR GENERATING A DATA
CLASSIFICATION MODEL USING AN ADAPTIVE LEARNING ALGORITHM**

Abstract of the Disclosure

5 A data classification method and apparatus are disclosed for labeling
unknown objects. The disclosed data classification system employs a learning algorithm
that adapts through experience. The present invention classifies objects in domain
datasets using data classification models having a corresponding bias and evaluates the
performance of the data classification. The performance values for each domain dataset
10 and corresponding model bias are processed to identify or modify one or more rules of
experience. The rules of experience are subsequently used to generate a model for data
classification. Each rule of experience specifies one or more characteristics for a domain
dataset and a corresponding bias that should be utilized for a data classification model if
the rule is satisfied. The present invention dynamically modifies the assumptions (bias)
15 of the learning algorithm to improve the assumptions embodied in the generated models
and thereby improve the quality of the data classification and regression systems that
employ such models. The disclosed self-adaptive learning process will become
increasingly more accurate as the rules of experience are accumulated over time.

20 1500-144.APP

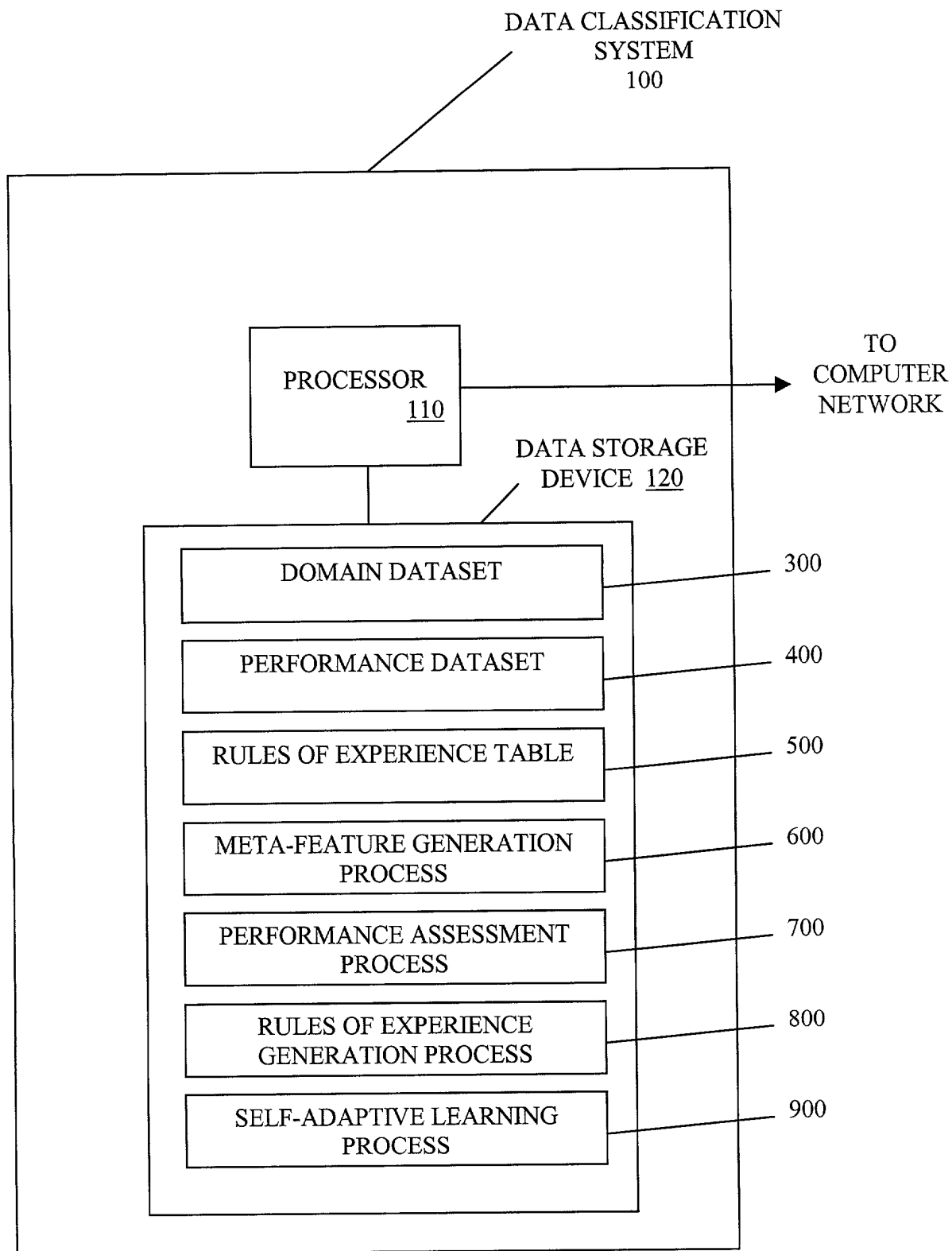


Figure 1

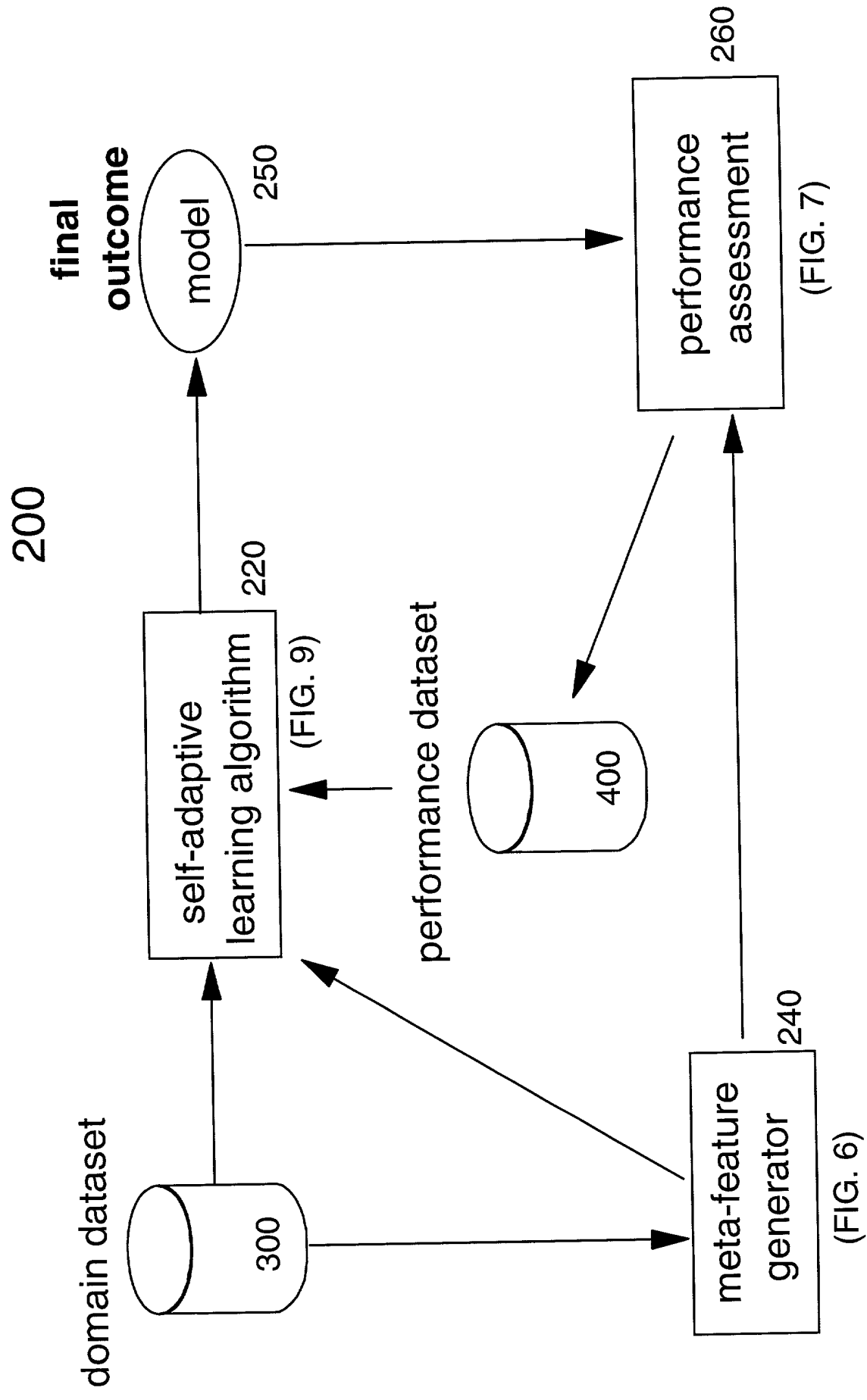


Figure 2

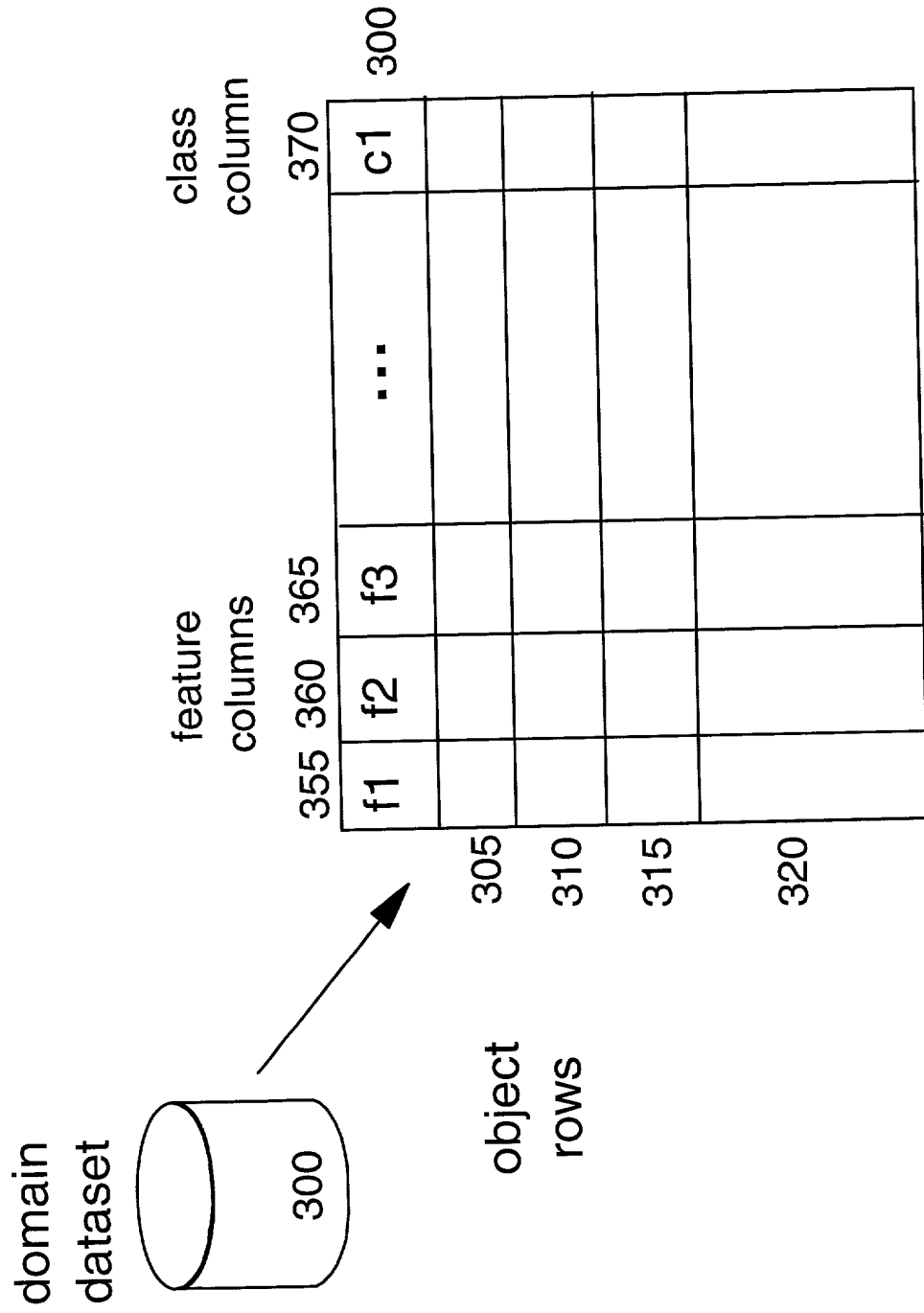


Figure 3

Performance Dataset		
450	455	460
405	meta-features domain 1	model bias 1
410	meta-features domain 2	model bias 2
415	meta-features domain 3	model bias 3
		quality 1
		quality 2
		quality 3
		400

Figure 4

550	560	570
Rule Identifier	Rule Condition	Bias
Rule 1	if density < 0.7 & variation > 0.8	non-linear discriminant
Rule 2	if entropy < 0.7 & variation < 0.4	linear discriminant
...		
Rule N	...	

Figure 5

600

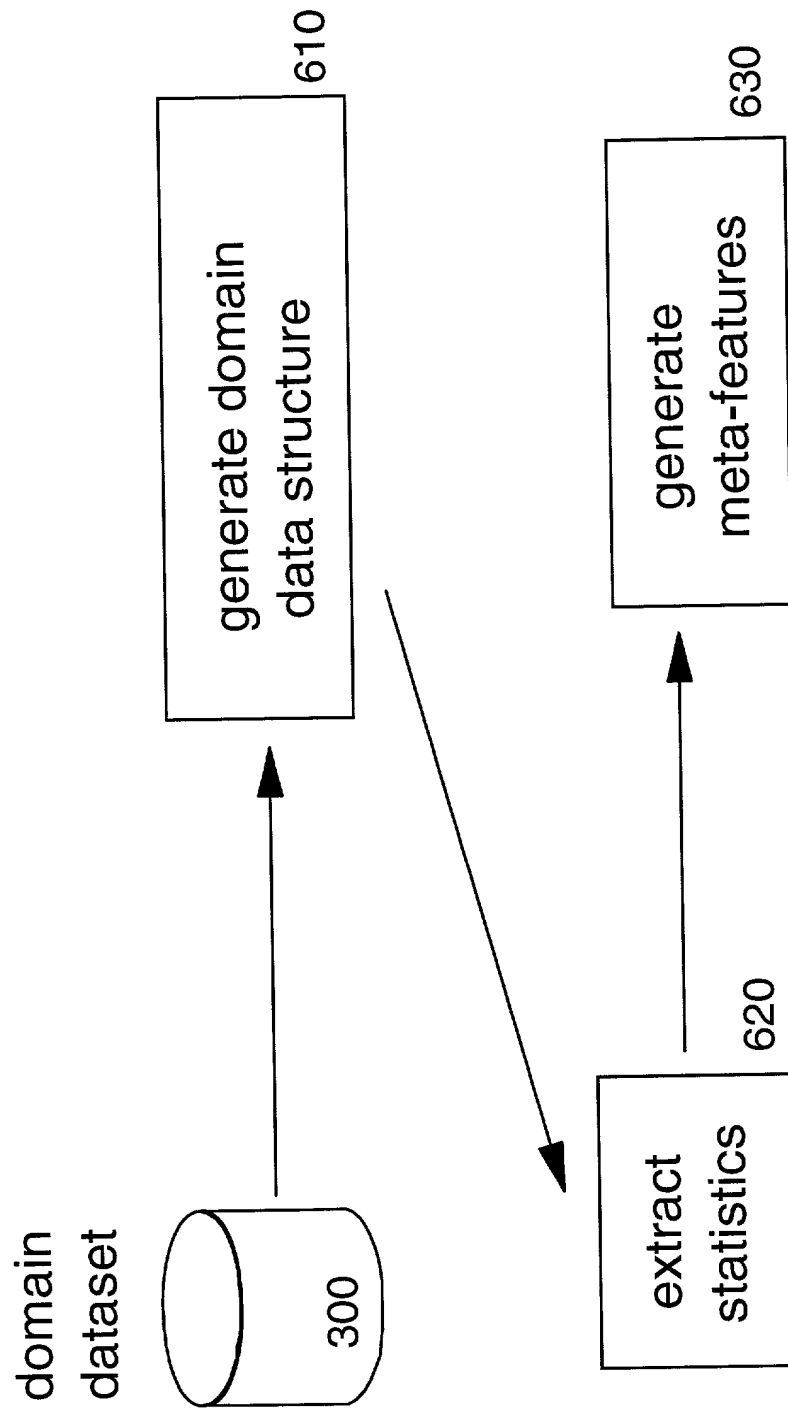
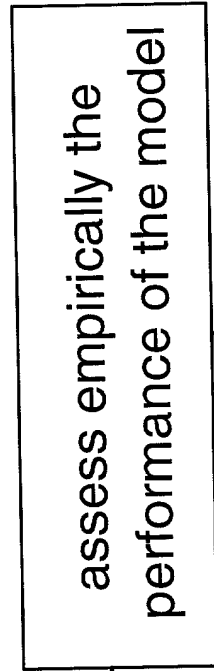


Figure 6

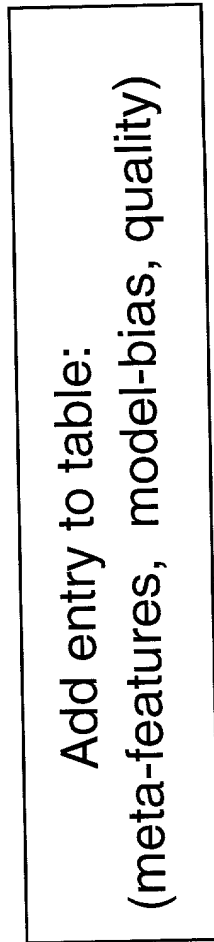
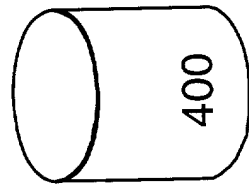
700

710



250

performance dataset



720

715

(Fig. 6)

Figure 7

800

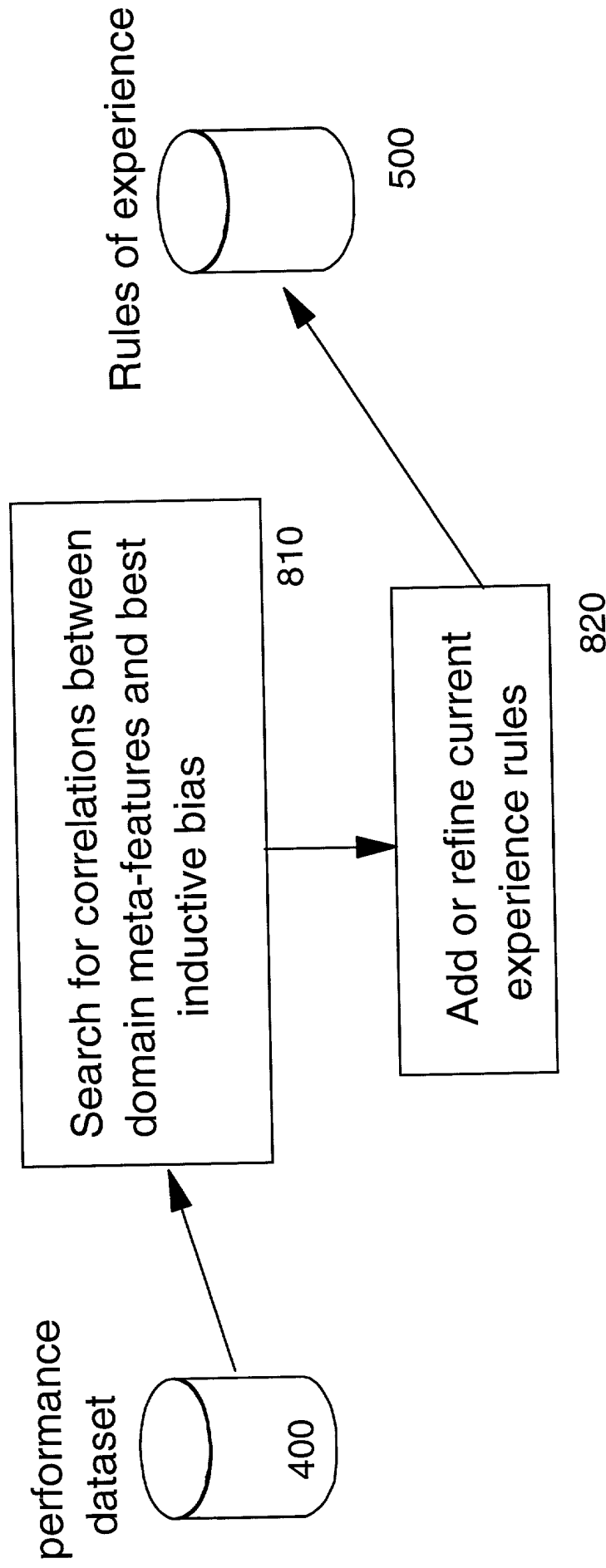


Figure 8

900

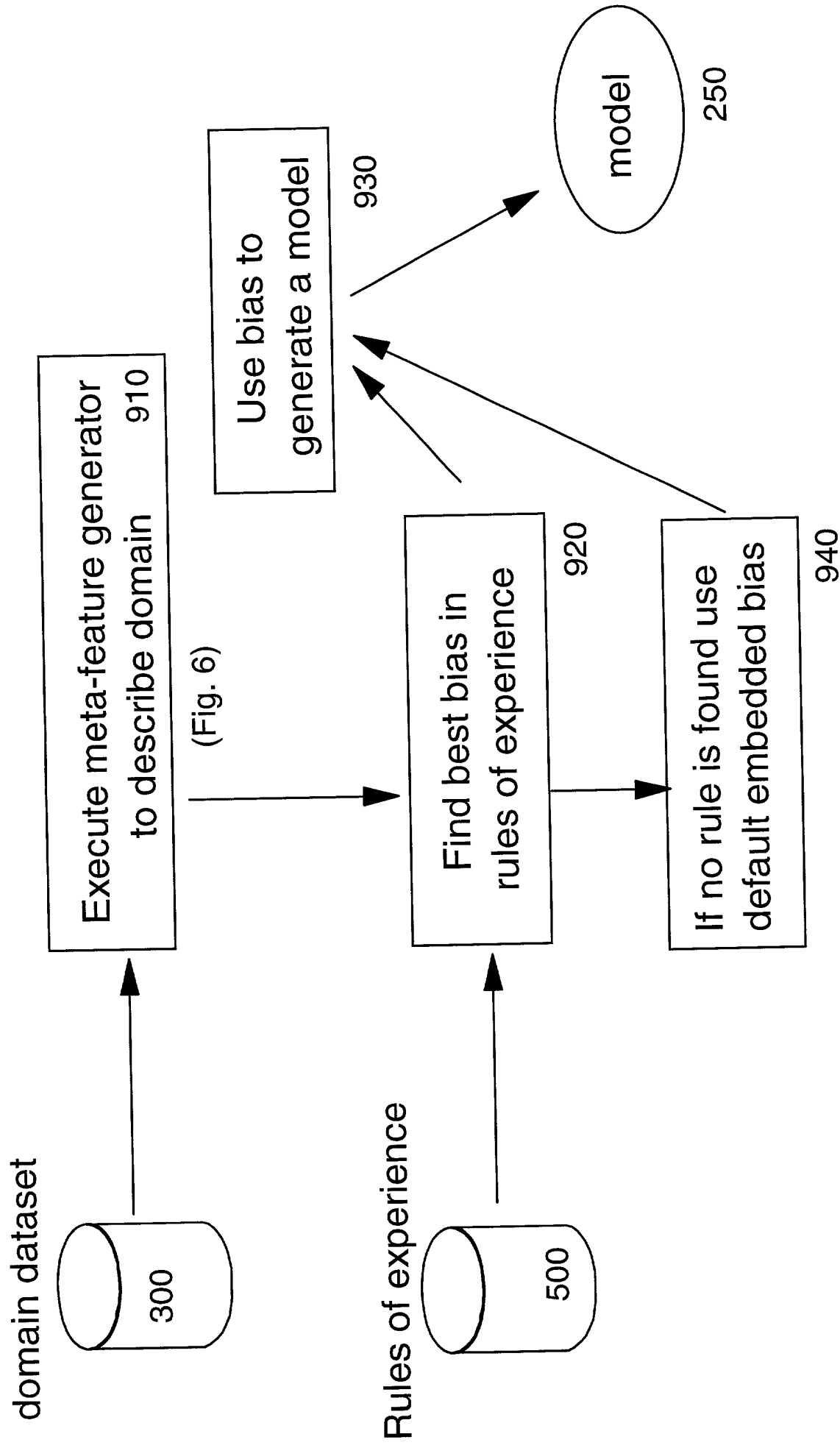


Figure 9

DECLARATION

AS A BELOW NAMED INVENTOR, I hereby declare that:

My residence, post office address and citizenship are as stated next to my name.

I believe that I am the original, first and sole (*if only one name is listed below*), or an original, first and joint inventor (*if plural names are listed below*), of the subject matter which is claimed and for which a patent is sought on the invention entitled:

TITLE: METHODS AND APPARATUS FOR GENERATING A DATA CLASSIFICATION MODEL USING AN ADAPTIVE LEARNING ALGORITHM

the specification of which is attached hereto or indicates an attorney docket no., or:

☐ was filed in the U.S. Patent & Trademark Office on _____ and assigned Serial No.,

☐ and (*if applicable*) was amended on _____.

I hereby state that I have reviewed and understand the contents of the above-identified specification, including the claims, as amended by any amendment referred to above. I acknowledge the duty to disclose information which is material to patentability and to the examination of this application in accordance with Title 37, Code of Federal Regulations §1.56. I hereby claim foreign priority benefits under Title 35, U.S. Code §119(a)-(d) or §365(b) of any foreign application(s) for patent or inventor's certificate, or §365(a) of any PCT international application which designated at least one country other than the United States, or §119(e) of any United States provisional application(s), listed below and have also identified below any foreign applications for patent or inventor's certificate having a filing date before that of the application on which priority is claimed:

Priority Claimed:

Yes [] No []

(Application Number) (Country) (Day/Month/Year filed)

Yes [] No []

(Application Number) (Country) (Day/Month/Year filed)


I hereby claim the benefit under Title 35, U.S. Code §120, of any United States application(s), or §365(c), of any PCT International application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States or PCT International applications(s) in the manner provided by the first paragraph of Title 35, U.S. Code §112, I acknowledge the duty to disclose information material to patentability as defined in Title 37, Code of Federal Regulations §1.56 which became available between the filing date of the prior application and the national or PCT international filing date of this application:

(Application Serial Number) (Filing Date) (STATUS: patented, pending, abandoned)

(Application Serial Number) (Filing Date) (STATUS: patented, pending, abandoned)

I hereby appoint the following attorneys: **MANNY W. SCHECTER**, Reg. No. 31,722; **LAUREN BRUZZONE**, Reg. No. 35,082; **CHRISTOPHER A. HUGHES**, Reg. No. 26,914; **EDWARD A. PENNINGTON**, Reg. No. 32,588; **JOHN E. HOEL**, Reg. No. 26,279; **JOSEPH C. REDMOND, Jr.**, Reg. No. 18,753; **DOUGLAS W. CAMERON**, Reg. No. 31,596; **LOUIS P. HERZBERG**, Reg. No. 41,500; **STEPHEN C. KAUFMAN**, Reg. No. 29,551; **DANIEL P. MORRIS**, Reg. No. 32,053; **PAUL J. OTTERSTEDT**, Reg. No. 37,411; **LOUIS J. PERCELLO**, Reg. No. 33,206; **ROBERT M. TREPP**, Reg. No. 25,933; and **MARIAN UNDERWEISER**, Reg. No. 46,134; each of them of **INTERNATIONAL BUSINESS MACHINES CORPORATION**, Thomas J. Watson Research Center, P.O. Box 218, Yorktown Heights, New York 10598; to prosecute this application and to transact all business in the U.S. Patent and Trademark Office connected therewith and with any divisional, continuation, continuation-in-part, reissue or re-examination application, with full power of appointment and with full power to substitute an associate attorney or agent, and to receive all patents which may issue thereon, and request that all correspondence be addressed to:

Kevin M. Mason
RYAN, MASON & LEWIS, LLP
1300 Post Road, Suite 205
Fairfield, CT 06430
Tel.: (203) 255-6560

FULL NAME OF FIRST OR SOLE INVENTOR: Youssef Drissi  Citizenship Morocco

Inventor's signature: Youngho Park Date: Nov. 18, 2022

Residence & Post Office address: 25 S. Highland Ave. #38
Ossining, NY 10562

FULL NAME OF SECOND JOINT INVENTOR: Ricardo Vilalta Citizenship Mexican

Inventor's signature: D. Z. [Signature] Date: Nov. 13, 2000

Residence & Post Office address: 255 Knickerbocker Ave.
Stamford, CT 06907

Attorney Docket No. YOR920000401US1

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICANT(S): Youssef Drissi and Ricardo Vilalta

SERIAL NO.: Unassigned

FILED: Concurrently Herewith

TITLE: METHODS AND APPARATUS FOR GENERATING A DATA
CLASSIFICATION MODEL USING AN ADAPTIVE LEARNING
ALGORITHM


ASSOCIATE POWER OF ATTORNEY

Please recognize **JOSEPH B. RYAN**, Reg. No. 37,922; **KEVIN M. MASON**, Reg. No. 36,597; **WILLIAM E. LEWIS**, Reg. No. 39,274; **ROBERT J. MAURI**, Reg. No. 41,180; and **WAYNE L. ELLENBOGEN**, Reg. No. 43,602; each of them of **RYAN, MASON & LEWIS, LLP**, 1300 Post Road, Suite 205, Fairfield, CT 06430 as associate attorneys in the above-mentioned application, with full power to prosecute said application, to make alterations and amendments therein, and to transact all business in the Patent and Trademark Office connected therewith.

Telephone calls should be made to Kevin M. Mason by dialing (203) 255-6560.

All written communications are to be sent to Kevin M. Mason, Esq., Ryan, Mason & Lewis, LLP, 1300 Post Road, Suite 205, Fairfield, CT 06430.

Dated: NOVEMBER 14, 2000


Paul J. Ottensmeyer
Registration No. 37,411
Attorney for Applicant(s)

International Business Machines Corporation
T.J. Watson Research Center
Route 134 and Kitchawan Road
Yorktown Heights, New York 10598

00477-2486